

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problem Mailbox.**

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平6-215184

(43) 公開日 平成6年(1994)8月5日

(51) Int.Cl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 K 9/36				
G 0 6 F 15/70	3 3 0 A	9071-5L		
G 0 6 K 9/20	3 4 0 K			
9/72		9289-5L		

審査請求 未請求 請求項の数 3 F D (全 10 頁)

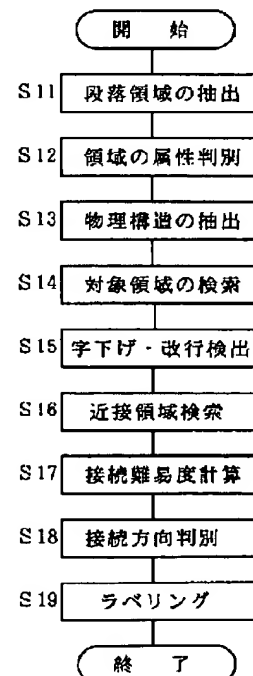
(21) 出願番号	特願平5-249966	(71) 出願人	000237156 富士ファコム制御株式会社 東京都日野市富士町1番地
(22) 出願日	平成5年(1993)9月10日	(71) 出願人	000005234 富士電機株式会社 神奈川県川崎市川崎区田辺新田1番1号
(31) 優先権主張番号	特願平4-273784	(72) 発明者	森 泰二 東京都日野市富士町1番地 富士ファコム 制御株式会社内
(32) 優先日	平4(1992)9月17日	(72) 発明者	本郷 保夫 東京都日野市富士町1番地 富士ファコム 制御株式会社内
(33) 優先権主張国	日本 (J P)	(74) 代理人	弁理士 森田 雄一

(54) 【発明の名称】 抽出領域のラベリング装置

(57) 【要約】

【目的】 文書画像から抽出した段落領域のラベリング精度を向上する。

【構成】 最初に、入力された文書画像から段落領域を抽出するとともに (S 1 1)、抽出された領域の属性を判別し (S 1 2)、属性が文字である段落領域についての物理構造を解析・抽出して探索木を作成する (S 1 3)。次に、属性が見出しである段落領域を挟む位置の文書段落領域を検索して対象領域とし (S 1 4)、抽出された各対象領域について、字下げ・改行の有無を検索する (S 1 5)。また、各段落領域が相互に接続する可能性のある近接領域を検索し (S 1 6)、さらに、各段落領域間の接続難易度をそれぞれ計算することにより (S 1 7)、対象の文書段落領域の接続方向を判別する (S 1 8)。得られた接続方向に基づき各文書段落領域の接続を行い、抽出領域をラベリングする (S 1 9)。



1

2

【特許請求の範囲】

【請求項1】 文書画像から抽出された段落領域の属性が文書であるか否かを判別する手段と、属性が文書であると判別された文書段落領域から探索木を作成する手段と、

属性が文書でないと判別された非文書段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせを探索木から検索する手段と、

検索された文書段落領域ごとに最終行の行末空白および先頭行の行頭空白を検出する手段と、

文書段落領域の先頭および末尾の空白の有無から接続の組合わせごとに文書段落領域間の接続の難易度を算出する手段と、

文書段落領域間の接続の組合わせごとの接続難易度を比較して非文書段落領域前後の文書段落領域の接続方向を判別する手段と、

判別された接続方向に基づき各文書段落領域のラベリングを行う手段と、

を備えたことを特徴とする抽出領域のラベリング装置。

【請求項2】 文書画像から抽出された段落領域の属性が文書であるか否かを判別する手段と、

属性が文書であると判別された文書段落領域から探索木を作成する手段と、

属性が文書でないと判別された非文書段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせを探索木から検索する手段と、

検索された文書段落領域ごとに文字認識を行う手段と、

文書段落領域の先頭および末尾の認識結果を接続の組み合わせごとに比較して両者の接続の適不適から文書段落領域間の接続の難易度を算出する手段と、

文書段落領域間の接続の組合わせごとの接続難易度を比較して非文書段落領域前後の文書段落領域の接続方向を判別する手段と、

判別された接続方向に基づき各文書段落領域のラベリングを行う手段と、

を備えたことを特徴とする抽出領域のラベリング装置。

【請求項3】 請求項1記載の抽出領域のラベリング装置において、請求項2記載の抽出領域のラベリング装置における文字認識手段および接続難易度算出手段を備え、文書段落領域間の接続の組合わせごとの接続難易度を接続方向別に比較し、接続難易度の値の差が所定値以下であれば、請求項2記載の文字認識手段および接続難易度算出手段により接続難易度を求めることを特徴とする抽出領域のラベリング装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、文字認識装置に入力された文書画像から文字領域を抽出した後、抽出領域をラベリングする抽出領域のラベリング装置に関する。

【0002】

10

20

30

40

50

【従来の技術】 従来の文字認識装置における文字領域の抽出およびラベリングは、図12に示すフローチャートのように行われている。図では、最初に文書画像が2値画像データとして入力されると(S21)、ぼかし等の処理を施した後、その周辺部のゴミ等を除去し(S22)、段組構造を判定する(S23)。次いで、縦書き・横書きを判定してから(S24)、行ブロックの抽出および行ブロックのラベリングを行う(S25、26)。さらに、段ブロックの抽出をした後に(S27)、段ブロックのラベリングを行う(S28)。この段ブロックのラベリングでは、段落の物理的な位置関係から探索木を作成し、予め設定されたルールに基づいて枝を検索して段落の接続順を決定していた。

【0003】

【発明が解決しようとする課題】 しかしながら、このような段落のラベリング方法では、新聞等の複雑な構造を有する文書に対して、必ずしも妥当なラベリングが行われるとは限らない。例えば、図13に示すような新聞記事の場合、見出し⑤が文字領域を区分するセパレータとして認識されてしまうと、その前後の段落③、⑥および④、⑦は見出し⑤を飛び越して接続されることがなく、誤ってラベリングされてしまう。同様に、図14に示すような新聞記事の場合、見出し⑥がセパレータとして認識されてしまうと、その前後の段落③、⑦および④、⑧は見出し⑥を飛び越して接続されることがなく、誤ってラベリングされてしまうという問題があった。本発明は上記問題点を解決するためになされたもので、その目的とするところは、見出しをセパレータとして認識することなく常に正確にラベリングすることのできる抽出領域のラベリング装置を提供することにある。

【0004】

【課題を解決するための手段】 上記目的を達成するために、第1の発明は、文書画像から抽出された段落領域の属性が文書であるか否かを判別する手段と、属性が文書であると判別された文書段落領域から探索木を作成する手段と、属性が文書でないと判別された非文書段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせを探索木から検索する手段と、検索された文書段落領域ごとに最終行の行末空白および先頭行の行頭空白を検出する手段と、文書段落領域の先頭および末尾の空白の有無から接続の組合わせごとに文書段落領域間の接続の難易度を算出する手段と、文書段落領域間の接続の組合わせごとの接続難易度を比較して非文書段落領域前後の文書段落領域の接続方向を判別する手段と、判別された接続方向に基づき各文書段落領域のラベリングを行う手段とを備えたことを特徴とする。

【0005】 第2の発明は、文書画像から抽出された段落領域の属性が文書であるか否かを判別する手段と、属性が文書であると判別された文書段落領域から探索木を作成する手段と、属性が文書でないと判別された非文書

段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせを探索木から検索する手段と、検索された文書段落領域ごとに文字認識を行う手段と、文書段落領域の先頭および末尾の認識結果を接続の組み合わせごとに比較して両者の接続の適不適から文書段落領域間の接続の難易度を算出する手段と、文書段落領域間の接続の組合わせごとの接続難易度を比較して非文書段落領域前後の文書段落領域の接続方向を判別する手段と、判別された接続方向に基づき各文書段落領域のラベリングを行う手段とを備えたことを特徴とする。

【0006】第3の発明は、第1の発明に、第2の発明における文字認識手段および接続難易度算出手段を備え、文書段落領域間の接続の組合わせごとの接続難易度を接続方向別に比較し、接続難易度の値の差が所定値以下であれば、第2の発明の文字認識手段および接続難易度算出手段により接続難易度を求めるようにしたことを特徴とする。

【0007】

【作用】第1の発明においては、文書画像から抽出された段落領域の属性が文書であるか否かが判別され、属性が文書であると判別された文書段落領域から探索木が作成される。属性が文書でないと判別された非文書段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせが探索木から検索される。次いで、検索された文書段落領域ごとに最終行の行末空白および先頭行の行頭空白が検出され、空白の有無から接続の組合わせごとに文書段落領域間の接続の難易度が算出される。さらに、算出された接続難易度が比較されて非文書段落領域前後の文書段落領域の接続方向が判別され、その接続方向に基づき各文書段落領域のラベリングが行われる。

【0008】第2の発明においては、文書画像から抽出された段落領域の属性が文書であるか否かが判別され、属性が文書であると判別された文書段落領域から探索木が作成される。属性が文書でないと判別された非文書段落領域の前後に位置しかつ他の文書段落領域との接続が可能な文書段落領域の組合わせが探索木から検索される。次いで、検索された文書段落領域ごとに文字認識が行われ、文書段落領域の先頭および末尾の認識結果が接続の組み合わせごとに比較され、両者の接続の適不適から文書段落領域間の接続の難易度が算出される。さらに、算出された接続難易度が比較されて非文書段落領域前後の文書段落領域の接続方向が判別され、その接続方向に基づき各文書段落領域のラベリングが行われる。

【0009】第3の発明においては、第1の発明に、第2の発明における文字認識手段および接続難易度算出手段が備えられており、文書段落領域間の接続の組合わせごとの接続難易度が接続方向別に比較され、接続難易度の値の差が所定値以下であれば、第2の発明の文字認識

手段および接続難易度算出手段により接続難易度が求められる。

【0010】

【実施例】以下、図に沿って本発明の実施例を説明する。図1は第1の発明の実施例の処理動作を示すフローチャートである。図において、最初に、入力された文書画像から段落領域を抽出するとともに(S11)、抽出された領域の属性を判別する(S12)。ここで判別される属性としては、文書、見出し、図、写真等がある。次いで、属性が文書である段落領域についての物理構造を解析・抽出する(S13)。ここで抽出された物理構造より探索木が作成される。

【0011】さらに、属性が文書ではなく見出し等である段落領域を挟む位置の文書段落領域を検索して対象領域とする(S14)。このように探索木を検索すると、例えば図2に示されるように、見出し領域1にかかる2段に、それぞれ前後して位置する文書段落領域2~5が抽出される。次に、抽出された各対象領域について、字下げ・改行の有無を検索する(S15)。具体的には、内蔵する行の座標や投影等の画像処理により、段落領域先頭の字下げと段落領域最後の改行による空白を検出してマークする。また、各段落領域が接続する可能性のある近接領域を検索する(S16)。この近接領域の検索は図3のように行われる。

【0012】図では、見出し領域6の前の段落領域7が見出し領域6を飛び越す場合に文書段落領域8へ接続され、飛び越さない場合に段落領域9へ接続される。同様に、見出し領域6の後の段落領域8の接続先についても、文書段落領域9へ接続される場合と文書段落領域10へ接続される場合がある。つまり、見出し領域6がセパレータとして機能する場合は、文書段落領域7~12等の接続方向がそれぞれ縦方向となるが、見出し領域6がセパレータとして機能しない場合の接続方向はそれぞれ横方向となる。

【0013】次に、各段落領域間の接続の組み合わせごとに接続難易度を計算する(S17)。具体的には図4に示すように、前の文書段落領域が改行で終了し後の文書段落領域が字下げで始まっている場合は接続の可能性が大であるから接続難易度を+1とし、前の文書段落領域が改行で終了し後の文書段落領域に字下げがない場合は接続の可能性が小であるから接続難易度を-1とする。また、前の文書段落領域に改行がなく、後の文書段落領域に字下げがある場合は接続の可能性が小であるから接続難易度を-1とする。

【0014】さらに、両文書段落領域とも空白がない場合は接続関係が不明であるものとして接続難易度を0とする。また、これ以外にも、両文書段落領域に空白がなくしかもその間に見出しが位置している場合は、両者が接続される可能性が極めて高いものとして接続難易度を1.5とする。その理由は、見出しがセパレータである

なら後の文書段落領域にも必ず字下げがあるはずだからである。さらに、上述した以外にも、文書段落領域の配置上のルールにより接続難易度を設定することが可能である。

【0015】次いで、得られた接続難易度を対象の文書段落領域の中の見出しの前または後における縦方向、および見出しをまたぐ横方向ごとにそれぞれ累積し、その縦と横の累積値を比較することにより対象の文書段落領域の接続方向を判別する(S18)。その結果、得られた接続方向に基づいて、各文書段落領域の接続を行い、抽出領域のラベリングを終了する(S19)。上述した実施例の方法で図13の文書段落領域をラベリングした結果が図5となる。この第1の発明の実施例では、文書段落領域の間に見出しがある場合でも、その前後の文書段落領域の空白の有無により接続方向が比較的高速に判定されて正しくラベリングが行われるようになる。

【0016】次に、第2の発明の実施例について説明する。図6は第2の発明の実施例の処理動作を示すフローチャートである。図において、最初に、入力された文書画像から段落領域を抽出するとともに(S31)、抽出された領域の属性を判別する(S32)。ここで判別される属性としては、文書、見出し、図、写真等がある。次いで、属性が文書である段落領域についての物理構造を解析・抽出して、探索木を作成する(S33)。図7は、図14に示す文書画像から作成された探索木を示す。図中のHは文書段落領域であるところの本文領域を、Mは見出し領域を表す。

【0017】さらに、属性が文書である領域について順次文字認識を行い(S34)、得られた文字コードを段落領域情報等とともにメモリに記憶する(S35)。次に、探索木より、図、写真、見出し等の本文以外の領域を検索し、着目領域とする(S36)。さらに着目領域が本文領域に挟まれているか否かによりセパレータであるか否かを判別し(S37)、挟まれていなければ着目領域がセパレータであるものとしてS43へ進みラベリングを行う(S38Yes)。

【0018】挟まれていれば(S38No)、着目領域に隣接する本文領域を探索木より検索し(S39)、さらに検索された本文領域相互間における接続の可能性を検索する(S40)。これらの近接領域の検索は第1の発明の実施例と同様に図3のように行われる。次に、検索により得られた本文領域間の接続難易度を求める(S41)。ここで求められる接続難易度とは、S34において領域ごとに認識された文章を言語処理技術を用いて単語に分割し、次いで、段落間をまたぐ単語について予め用意した単語辞書との照合結果および文法ルールによる文法チェックの結果を接続確率等のように数値化し接続難易度としたものである。図8は、これらの接続難易度計算の処理順を示したフローチャートである。

【0019】また、他にも本文領域末尾の文字の文字種

(英数、カタカナ)が、対応する本文領域先頭の文字の文字種と一致しているか否かによっても接続難易度が得られる。さらには、括弧が開いたまま終わっている領域は、閉じ括弧から始まる領域に接続する可能性が大きいように、括弧記号がある場合の位置関係からも接続難易度が得られる。これらの方法により、図14の文書画像について、見出し領域の前後の文書領域について接続難易度を求めると、図9のようになる。これらの結果から、図14の文書画像について文書領域間の接続難易度を図示すると、図10のようになる。図中の実線は難易度+1を、破線は難易度0を表す。

【0020】次に、得られた接続難易度を、着目領域をまたぐ方向(横方向)と着目領域に従う方向(縦方向)にそれぞれ累積し、得られた累積値を比較することで接続方向を判別する(S42)。すなわち、着目領域をまたぐ方向の接続難易度の累積値が大きければ、着目領域はセパレータとして機能しないものと判別される。次に得られた接続方向に基づき、本文領域間の接続がなされ抽出領域のラベリングを終了する(S43)。これらの処理により図14の文書画像をラベリングした結果が図11となる。

【0021】この第2の発明の実施例では、文書段落領域の間に見出しがある場合でも、その前後の文書段落領域の文字を認識して接続難易度を算出することにより、接続方向が確実に判定されて正しくラベリングが行われるようになる。なお、第2の発明の実施例では、言語処理の対象とする部分を段落領域の比較対象となる先頭と末尾に限定することで処理の高速化が可能になる。さらには、図6のフローチャートにおけるS33とS34、S34とS36～S40は並列処理が可能であるので処理速度をさらに高速にすることも可能である。また、S34では全ての本文領域について文字認識をしているが、S40で接続の可能性ありとされた本文領域についてのみ文字認識をすればさらに高速化が可能である。

【0022】次に、第3の発明の実施例について説明する。前述した第2の発明の実施例では、接続難易度を算出する際に単語辞書を参照するためにメモリのアクセス時間が長くなり、高速化に限界がある。そこで第3の発明では、処理速度を増すため、第1の発明および第2の発明でそれぞれ求めた接続難易度を用いて接続方向を判定してラベリングするようにした。つまり、通常は第1の発明により比較的高速で接続難易度を求めて接続方向を判定するが、第1の発明の接続難易度だけでは十分に判定できない場合に第2の発明で求めた接続難易度により接続方向を判定してラベリングするようにしたものである。それにより、確実に高速なラベリングが可能になる。

【0023】

【発明の効果】以上述べたように第1の発明によれば、見出し段落領域前後に位置する文書段落領域の先頭およ

び末尾の空白の有無から文書段落領域間の接続の難易度が算出されて、見出し段落領域前後の文書段落領域の接続方向が判別される。それにより、新聞等の構造が複雑な文書であっても、見出しをセパレータとして誤認識することなく高速にラベリングすることができる。

【0024】第2の発明によれば、見出し段落領域前後に位置する文書段落領域の先頭および末尾の認識結果を互いに比較して文書段落領域間の接続の難易度を算出することにより、見出し段落領域前後の文書段落領域の接続方向が判別される。それにより、新聞等の構造が複雑な文書であっても、見出しをセパレータとして誤認識することなく正確にラベリングすることができる。

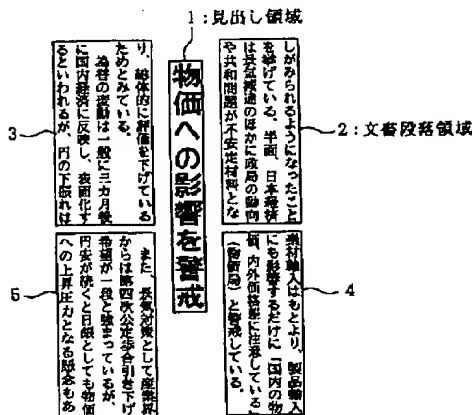
【0025】第3の発明によれば、第1の発明と第2の発明を組み合わせ、通常は比較的処理速度の速い第1の発明により文書段落領域間の接続難易度を求めて接続方向を判別する。さらに、第1の発明では方向ごとの接続難易度の値に明瞭な差が認められない場合に、第2の発明を用いて文書段落領域の文字を認識して接続難易度を再度求めることにより高精度に接続難易度を求めて接続方向を判別する。その結果、入力文書画像から抽出された文書領域を高速かつ高精度にラベリングすることが可能になる。

【図面の簡単な説明】

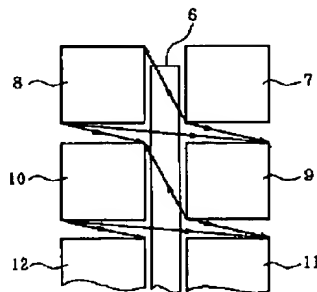
【図1】第1の発明の実施例の処理動作を示すフローチャートである。

【図2】第1の発明の実施例におけるラベリング対象の段落領域の説明図である。

【図2】



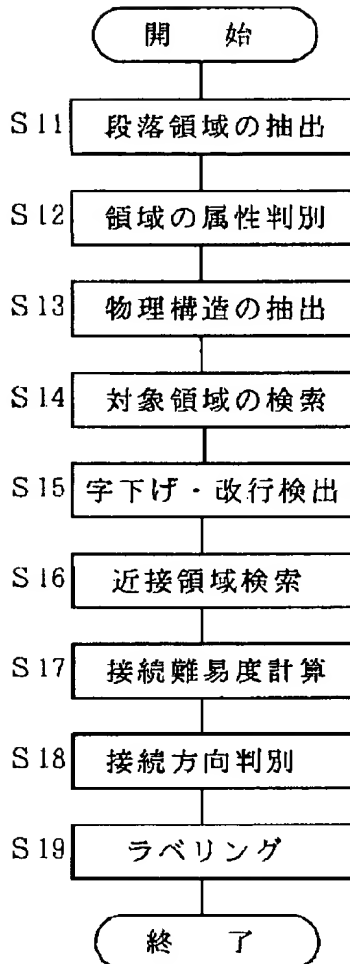
【図3】



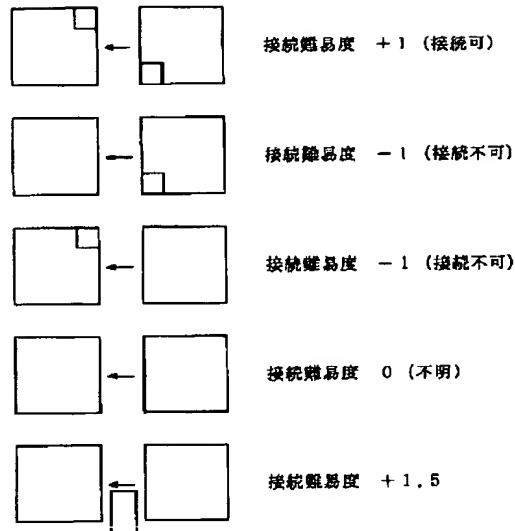
【図9】

		接続難易度
八割を占	占める	+1
	軍の指揮	0
和解、実	藤は	+1
	閉鎖した	0
<u>言葉処理の例</u>		
ソフト	ウェア	+1
	特許	0
UN	TAC	+1
	会議	0
<u>文字種の例</u>		
「部分 → 和平」		+1
「ポト派 → 部分和平」		-1
<u>括弧の位置の例</u>		

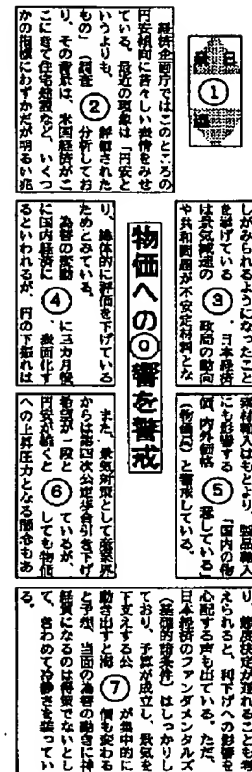
【図1】



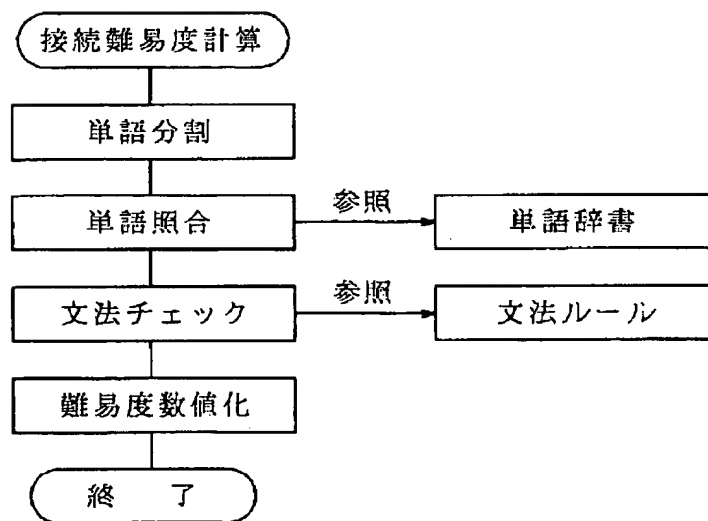
【図4】



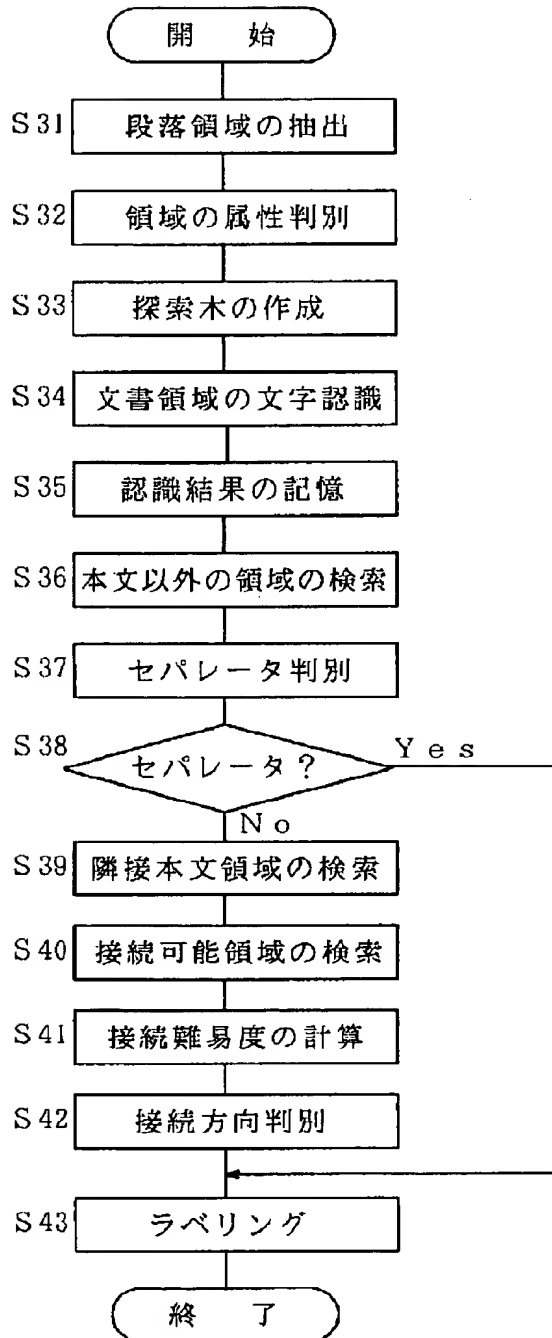
【図5】



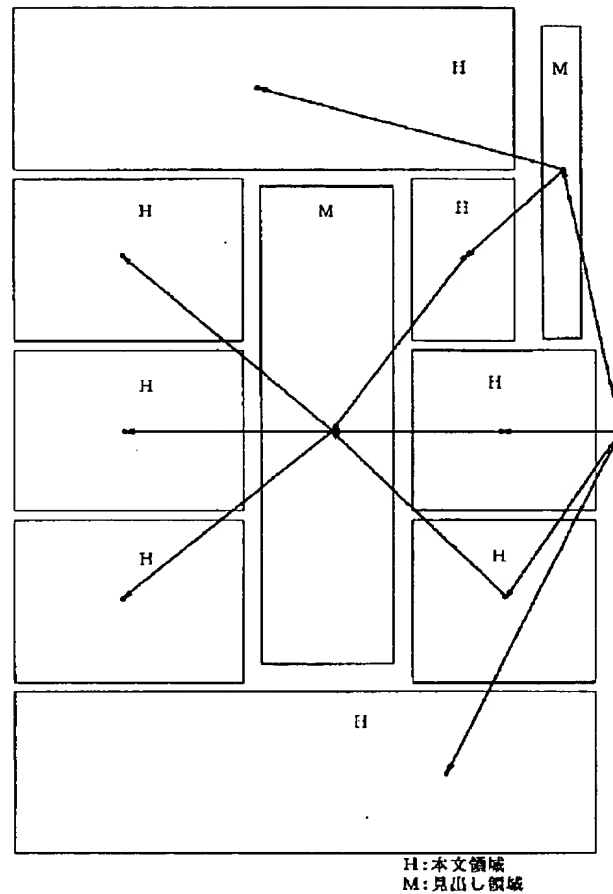
【図8】



【図6】



【図7】



【图 1 1】

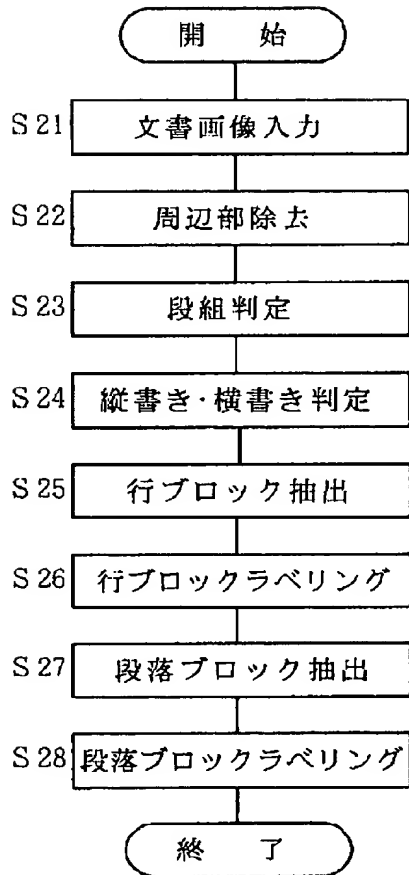
カンボジア①PKO

[illegible]

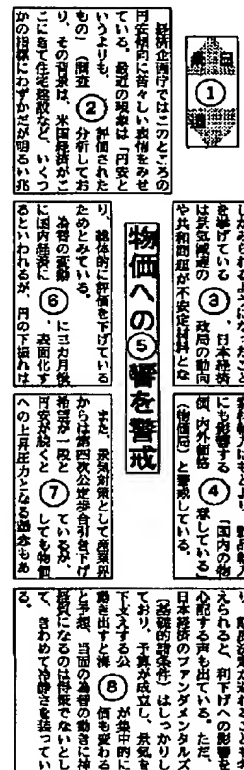
難しい日本の対応
選挙妨害で緊迫

[illegible][illegible]

【図12】



【図13】



【図14】

カンボジアのPKO

難しい日本の対応 選挙妨害で緊迫

カンボジアは内戦が入り、約100万人の難民が、国境を越えて隣国に逃げた。PKOに参入する日本は、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。カンボジアは、内戦が激化し、約100万人の難民が、国境を越えて隣国に逃げた。PKOに参入する日本は、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

① 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

② 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

③ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

④ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑤ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑥ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑦ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑧ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑨ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。

⑩ 日本は、カンボジアに、二千人の兵力を、二カ国に展開する。この二カ国は、カンボジアとラオスである。